

# OPTIMIZATION WITH TENSORS

General problem: Given a loss function  $\mathcal{L}: \mathbb{R}^{d_1 \times \dots \times d_N} \rightarrow \mathbb{R}$ , we want to solve

$$\min_{W \in \mathbb{R}^{d_1 \times \dots \times d_N}} \mathcal{L}(W) \quad \text{subject to} \quad \text{rank}(W) \leq R$$

↳ CP rank, Tucker rank, TT rank, ...

Examples of loss functions:

+ Low rank approximation

$$\min_W \|W - T\|_F^2 \quad \text{s.t.} \quad \text{rank}_{\text{CP}}(W) \leq R$$

↳ the loss function is  $\mathcal{L}(W) = \|W - T\|_F^2$

+ Regression

We want to learn a linear function  $f: \mathbb{R}^{d_1 \times \dots \times d_N} \rightarrow \mathbb{R}$

from a dataset  $\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\} \subseteq \mathbb{R}^{d_1 \times \dots \times d_N} \times \mathbb{R}$ .

$$f \text{ is linear} \Leftrightarrow f(x) = f_W(x) = \langle W, x \rangle = \sum_{i_1, \dots, i_N} W_{i_1, \dots, i_N} x_{i_1, \dots, i_N}.$$

In this case, a natural loss function is

$$\begin{aligned} \mathcal{L}(W) &= \sum_{i=1}^m (f_W(x_i) - y_i)^2 \\ &= \sum_{i=1}^m (\langle W, x_i \rangle - y_i)^2 \end{aligned}$$

+ Completion:

Target tensor  $T \in \mathbb{R}^{d_1 \times \dots \times d_N}$

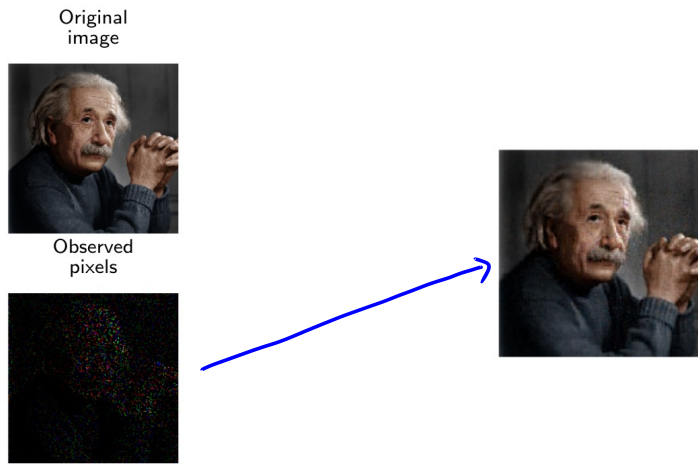
Data: observed entries  $\{T_{i_1, \dots, i_N} \mid (i_1, \dots, i_N) \in \Omega\}$

$$\Omega \subseteq [d_1] \times [d_2] \times \dots \times [d_N]$$

$$\bullet \begin{pmatrix} 1 & ? & ? \\ ? & 2 & ? \\ ? & ? & 3 \end{pmatrix}$$

$$\Omega = \{(1,1), (2,2), (2,3)\}$$

- Image completion



For completion, the loss function is:

$$\mathcal{L}(W) = \sum_{(i_1, \dots, i_n) \in \Omega} (T_{i_1, \dots, i_n} - W_{i_1, \dots, i_n})^2$$

↳ set of observed entries

Remark: If  $\Omega = [d_1] \times [d_2] \times \dots \times [d_n]$ , then

$$\sum_{(i_1, \dots, i_n) \in \Omega} (T_{i_1, \dots, i_n} - W_{i_1, \dots, i_n})^2 = \|T - W\|_F^2$$

- COLLABORATIVE FILTER / RECOMMENDATION SYSTEMS

(NETFLIX CHALLENGE)



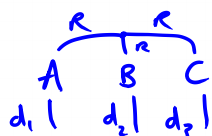
user i gave rating 4 on movie j

# ① GENERAL OPTIMIZATION ALGORITHM

**PBT 1**  $\min_{W \in \mathbb{R}^{d_1 \times \dots \times d_N}} \mathcal{L}(W)$  subject to  $\text{rank}_{CP}(W) \leq R$

Notation: Given  $A_i \in \mathbb{R}^{d_i \times R}$  for  $i=1, \dots, N$ , we denote the CP decomposition with factors  $A_1, \dots, A_N$  by  $CP(A_1, \dots, A_N)$

ex:  $CP(A, B, C) =$



PBT 1 is equivalent to:

**PBT 2**  $\min_{\substack{A_i \in \mathbb{R}^{d_i \times R} \\ i=1, \dots, N}} \mathcal{L}(CP(A_1, \dots, A_N))$

## 1) Gradient based algorithm

**ALGORITHM**: Initialize  $A_1, A_2, \dots, A_N$

Repeat  
  for  $i=1 \dots N$   
     $X_i \leftarrow A_i - \gamma \nabla_{A_i} \mathcal{L}(CP(A_1, \dots, A_N))$  ← learning rate  
  for  $i=1 \dots N$   
     $A_i \leftarrow X_i$   
  until convergence

## 2) Alternating minimization

We assume  $\min_{A_i} \mathcal{L}(CP(A_1, \dots, A_N))$  is easy.

**ALGORITHM**: Initialize  $A_1, A_2, \dots, A_N$

Repeat  
  for  $i=1, \dots, N$   
     $A_i \leftarrow \arg \min_{A_i} \mathcal{L}(CP(A_1, \dots, A_N))$   
  until convergence

### 3) Alternating minimization for low CP rank approximation

Given  $T \in \mathbb{R}^{d_1 \times \dots \times d_N}$ , target rank  $R$ :

**CP-PBM**

$$\min_{\substack{A_i \in \mathbb{R}^{d_i \times R} \\ i=1, \dots, N}} \|T - CP(A_1, \dots, A_N)\|_F^2$$

• Matricization:  $T \in \mathbb{R}^{d_1 \times \dots \times d_N}$ , for any  $m \in [N]$ ,  $T_{(m)} \in \mathbb{R}^{d_m \times d_1 \dots d_{m-1} d_{m+1} \dots d_N}$   
 "mode- $m$  matricization" ↗

• Remark:  $\|T\|_F = \|T_{(m)}\|_F$  for all  $m \in [N]$

Def: If  $A \in \mathbb{R}^{m \times m}$  and  $B \in \mathbb{R}^{n \times q}$ , their Kronecker product  
 $A \otimes B \in \mathbb{R}^{mn \times nq}$  is defined by:

$$A \otimes B = \begin{pmatrix} a_{11}B & \dots & a_{1m}B \\ \vdots & \ddots & \vdots \\ a_{m1}B & \dots & a_{mn}B \end{pmatrix}$$

$mn$   $nq$

( $a_{ij}$  is the entry  $(i,j)$  of  $A$ )

Remark:  $\underbrace{a}_{m \times 1} \otimes \underbrace{b}_{n \times 1} = \text{vec}(b a^T) = \text{vec}(\underbrace{b \circ a}_{m \times n})$   
 $mn \times 1$

Def: The Kathri-Rao product of  $A = \begin{pmatrix} | & & | \\ a_1 & \dots & a_R \\ | & & | \end{pmatrix} \in \mathbb{R}^{m \times R}$  and  
 $B = \begin{pmatrix} | & & | \\ b_1 & \dots & b_R \\ | & & | \end{pmatrix} \in \mathbb{R}^{m \times R}$  is defined by

$$A \odot B = \begin{pmatrix} | & & | \\ a_1 \otimes b_1 & a_2 \otimes b_2 & \dots & a_R \otimes b_R \\ | & & | \end{pmatrix} \in \mathbb{R}^{mm \times R}$$

Property: If  $W = CP(A_1, A_2, \dots, A_N)$ , then

$$W_{(m)} = \underbrace{A_m}_{d_m \times d_N \dots d_{m+1} d_{m-1} \dots d_1} \left( \underbrace{A_N \odot \dots \odot A_{m+1} \odot A_{m-1} \odot \dots \odot A_1}_{d_N \dots d_{m+1} d_{m-1} \dots d_1 \times R} \right)^T$$

## Alternating least squares algorithm (ALS) for CP decomposition

$$\min_{\substack{A, B, C \\ d_1 \times R \quad d_2 \times R \quad d_3 \times R}} \|T - CP(A, B, C)\|_F^2$$

We want solve this problem w.r.t  $A$ :

$$\begin{aligned} \arg\min_A \|T - CP(A, \overset{\text{fixed}}{B, C})\|_F^2 &= \arg\min_A \|T_{(1)} - (CP(A, B, C))_{(1)}\|_F^2 \\ &= \arg\min_A \|T_{(1)} - A(C \odot B)^T\|_F^2 \\ &= T_{(1)} [(C \odot B)^T]^+ \end{aligned}$$

pseudo-inverse

(side note: if  $A \in \mathbb{R}^{m \times n}$  with  $m \leq n$  and  $\text{rank}(A) = m$ , then  $AA^+ = I$  (however:  $\triangle A^+A$  may not be  $I$   $\triangle$ ))

### ALS ALGORITHM:

INPUT: target tensor  $T$ , rank  $R$

OUTPUT: factor matrices  $A_i, i=1, \dots, N$  s.t.  
 $T \approx CP(A_1, \dots, A_N)$

- Initialize  $A_1, \dots, A_N$

- Repeat

$$\left[ \begin{array}{l} \text{For } m = 1, \dots, N \\ \quad A_m \leftarrow T_{(m)} \left[ (A_N \odot \dots \odot A_{m+1} \odot A_{m-1} \odot \dots \odot A_1)^T \right]^+ \\ \text{until convergence} \end{array} \right.$$

## II SVD-based algorithms for TUCKER and TT

### 1) TUCKER: Higher order SVD (LAUTHAUER et al. 2000)

Given  $T \in \mathbb{R}^{d_1 \times \dots \times d_N}$  and  $(R_1, \dots, R_N)$  we want to solve

$$\min_{\substack{G \in \mathbb{R}^{R_1 \times \dots \times R_N} \\ U_i \in \mathbb{R}^{d_i \times R_i}, \\ i=1 \dots N}} \| T - \underbrace{G \times_1 U_1 \times_2 U_2 \times \dots \times_N U_N}_{\text{Diagram}} \|_F^2$$

A diagram showing a central node 'G' with arrows pointing to 'U1', 'U2', and 'UN'.

### HOSVD Algorithm

INPUT:  $T \in \mathbb{R}^{d_1 \times \dots \times d_N}$  and  $(R_1, \dots, R_N)$

OUTPUT:  $G \in \mathbb{R}^{R_1 \times \dots \times R_N}$  s.t.  $T \approx G \times_1 U_1 \times_2 \dots \times_N U_N$   
 $U_i \in \mathbb{R}^{d_i \times R_i},$   
 $i=1 \dots N$

For  $m=1 \dots N$

$\begin{cases} U_m \leftarrow R_m \text{ leading left singular vectors of } T_{(m)} \\ d_m \times R_m \end{cases}$

$$G \leftarrow \underbrace{T}_{d_1 \times \dots \times d_N} \times_1 \underbrace{U_1^T}_{R_1 \times d_1} \times_2 \dots \times_N \underbrace{U_N^T}_{R_N \times d_N}$$

$R_1 \times \dots \times R_N \rightarrow$

RETURN  $G, U_1, \dots, U_N$

$\downarrow$  SVD  
 $U D V^T$   
 $\downarrow$   
 extract the first  $R_m$  columns

$$\star \min_X \| T - X \|_F^2 \quad \text{s.t.} \quad \text{rank}_{\text{TUCKER}}(X) \leq (R_1, \dots, R_N)$$

Theorem: Let  $X^*$  be the solution of problem  $\star$  and  $X_{\text{HOSVD}}$  the solution returned by HOSVD.  
 $\hookrightarrow G \times_1 U_1 \times \dots \times U_N$

Then

$$\| T - X_{\text{HOSVD}} \|_F \leq \sqrt{N} \| T - X^* \|_F$$

• If  $T$  has Tucker rank less than  $(R_1, \dots, R_N)$ , then  $\| T - X^* \|_F = 0$   
 hence  $\| T - X_{\text{HOSVD}} \|_F = 0$   $\Rightarrow$  HOSVD recovers the exact TUCKER decomposition if it exists.

2) TENSOR TRAIN: TT-SVD (OSSEDELETS, 2010)  
 ↳ (a.k.a. MATRIX PRODUCT STATE (MPS): sequential SVD)

$$\min_{G_1, G_2, \dots, G_N} \|T - \text{TT}(G_1, G_2, \dots, G_N)\|_F^2$$

$$\hookrightarrow \underset{d_1}{G_1} \underset{d_2}{\overset{R_1}{-}} \underset{d_2}{G_2} \underset{d_3}{\overset{R_2}{-}} \underset{d_3}{G_3} \underset{d_4}{\overset{R_3}{-}} \dots \underset{d_N}{\overset{R_{N-1}}{-}} G_N$$

TT-SVD Algorithm

INPUT:  $T \in \mathbb{R}^{d_1 \times d_2 \times d_3 \times d_4}$ ,  $(R_1, \dots, R_3)$

OUTPUT:  $G_1, \dots, G_4$  s.t.  $T \approx \text{TT}(G_1, G_2, G_3, G_4)$

$$\bullet \underset{d_4}{\overset{d_1}{T}} \underset{d_3}{\overset{d_2}{-}} \approx \underset{d_1}{G_1} \underset{d_2}{\overset{R_1}{-}} \underset{d_4}{A} \underset{d_3}{\overset{d_2}{-}} \quad (\text{Low rank matrix approx.} \rightarrow \text{SVD})$$

$$\bullet \underset{R_1}{\overset{d_2}{A}} \underset{d_4}{\overset{d_3}{-}} \approx \underset{R_1}{\overset{d_2}{G_2}} \underset{R_2}{\overset{R_1}{-}} \underset{d_4}{B} \underset{d_3}{\overset{d_2}{-}}$$

$$\bullet \underset{R_2}{\overset{d_3}{B}} \underset{d_4}{\overset{d_3}{-}} \approx \underset{R_2}{\overset{d_3}{G_3}} \underset{R_3}{\overset{R_2}{-}} \underset{d_4}{G_4} \underset{d_3}{\overset{d_2}{-}}$$

• RETURN  $G_1, G_2, G_3, G_4$

Theorem: If  $X^*$  is the best approximation of  $T$  of TT rank  $(R_1, \dots, R_{N-1})$  then  $\|T - X_{\text{TT-SVD}}\|_F \leq \sqrt{N-1} \|T - X^*\|_F$ .

↳ solution returned by TT-SVD

$$X_{\text{TT-SVD}} = \underset{d_1}{G_1} \underset{d_2}{\overset{R_1}{-}} \underset{d_2}{G_2} \underset{d_3}{\overset{R_2}{-}} \dots \underset{d_N}{\overset{R_{N-1}}{-}} \underset{d_N}{G_N}$$

• If  $\|T - X^*\|_F = 0$  then TT-SVD returns an exact TT decomposition of  $T$ .

# III SUPERVISED LEARNING WITH TENSOR NETWORKS

Supervised Learning With Quantum-Inspired Tensor Networks

(NeurIPS, 2016)

E. Miles Stoudenmire<sup>1,2</sup> and David J. Schwab<sup>3</sup>

We want to learn a linear function  $f: \mathbb{R}^{d^m} \rightarrow \mathbb{R}^n$   
 input space is very high dimension

$$f: x \mapsto Wx$$

$\uparrow$   
 $n \times d^m$

Feature map: Image  $\rightarrow \mathbb{R}^{d \times d \times \dots \times d}$

$$x \in \mathbb{R}^{h \times w} \mapsto \underbrace{\phi(x_{1,1}) \circ \phi(x_{1,2}) \circ \phi(x_{1,3}) \circ \dots \circ \phi(x_{h,w})}_{\text{tensor of order } hw}$$

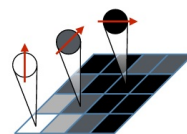
$\rightarrow \mathbb{R}^{2^{hw}}$

Two examples for  $\phi$  are

$$\begin{aligned} \phi: \mathbb{R} &\rightarrow \mathbb{R}^2 \\ \alpha &\mapsto \begin{pmatrix} \alpha \\ 1 \end{pmatrix} \end{aligned}$$

$$\hookrightarrow (\alpha_1, \alpha_2, \alpha_3) \mapsto \underbrace{\phi(\alpha_1) \circ \phi(\alpha_2) \circ \phi(\alpha_3)}_{\substack{\text{Entries of this tensor: } 1, \alpha_1, \alpha_2, \alpha_3, \alpha_1\alpha_2, \alpha_1\alpha_3, \alpha_2\alpha_3, \\ \text{and } \alpha_1\alpha_2\alpha_3}}$$

$$\begin{aligned} \phi: \mathbb{R} &\rightarrow \mathbb{R}^2 \\ \alpha &\mapsto \begin{pmatrix} \cos\left(\frac{\pi}{2}\alpha\right) \\ \sin\left(\frac{\pi}{2}\alpha\right) \end{pmatrix} \end{aligned}$$





$$f: \mathbb{R}^{d^m} \approx \mathbb{R}^{\underbrace{d \times d \times \dots \times d}_m} \longrightarrow \mathbb{R}^n$$

$$\underbrace{d \times d \times \dots \times d}_m \xrightarrow{\quad} \underbrace{d \times d \times \dots \times d}_m$$

↳ The model parameter is  $W \in \mathbb{R}^{d \times \dots \times d \times n}$

↳  $d^m n$  parameters!

We cannot even store  $W$  in memory!

Sol.: Parameterize  $W$  with a low rank TT decomposition

$$\underbrace{d \times d \times \dots \times d}_m \times n = \underbrace{d \times d \times \dots \times d}_m \times n$$

$$f(x) = f(\phi(x_{i1}) \circ \phi(x_{i2}) \circ \dots \circ \phi(x_{i n_w}))$$

$$= \underbrace{\phi(x_{i1}) \quad \phi(x_{i2}) \quad \dots \quad \phi(x_{i n_w})}_{\substack{d \times d \times \dots \times d \\ W \\ n}}$$

$$= \underbrace{\phi(x_{i1}) \quad \phi(x_{i2}) \quad \dots \quad \phi(x_{i n_w})}_{\substack{d \times d \times \dots \times d \\ G_1 \quad G_2 \quad \dots \quad G_{n_w} \quad G_{n_w+1} \\ n}}$$

← We can compute this very efficiently

Learning problem: Given data  $\{(x_i, y_i), \dots, (x_m, y_m)\} \subseteq \mathbb{R}^{d^m} \times \mathbb{R}^n$

we want to minimize  $\mathcal{L} = \sum_{i=1}^m \ell(f(x_i), y_i)$  for some loss function

$\ell: \mathbb{R}^n \rightarrow \mathbb{R}$ .

↳ parameters:  $G_1, G_2, \dots, G_{m+1}$  cores of a TT decomposition of  $W$ .

# Training Algorithm (DMRG)

Initialize  $G_1, G_2, \dots, G_{m+1}$

Repeat

for each consecutive pair of cores  $G_i, G_{i+1}$

(merge two cores)

$$\underbrace{R_{i-1}}_d \underbrace{\beta}_{d \times d} \underbrace{R_{i+1}}_d = \underbrace{R_{i-1}}_d \underbrace{G_i}_{d \times d} \underbrace{G_{i+1}}_d \underbrace{R_{i+1}}_d$$

(Gradient descent step)

$$\beta^{\text{new}} = \beta - \gamma \nabla_{\beta} \mathcal{L}$$

(split in two cores)

$$\underbrace{\beta^{\text{new}}}_{d \times d} \approx \underbrace{R_{i-1}}_d \underbrace{G_i^{\text{new}}}_{d \times d} \underbrace{R_i^{\text{new}}}_{d \times d} \underbrace{G_{i+1}^{\text{new}}}_{d \times d} \underbrace{R_{i+1}}_d \quad (\text{truncated SVD})$$

$$G_i, G_{i+1} \leftarrow G_i^{\text{new}}, G_{i+1}^{\text{new}}$$

until convergence

\*  $R_i^{\text{new}}$  can be chosen adaptively from the singular values of  $\beta^{\text{new}}$ .

Follow up papers:

From probabilistic graphical models to generalized tensor networks for supervised learning

Ivan Glasser,<sup>1,2</sup> Nicola Pancotti,<sup>1,2</sup> and J. Ignacio Cirac<sup>1,2</sup>

## Tensor Networks for Probabilistic Sequence Modeling

Jacob Miller  
Mila and DIRO  
Université de Montréal  
Montréal QC, Canada  
jmjacobmiller@gmail.com

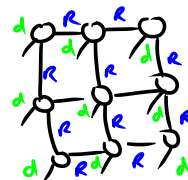
Guillaume Rabusseau  
Mila and DIRO  
Université de Montréal  
Montréal QC, Canada  
grabus@iro.umontreal.ca

John Terilla  
CUNY and TUNNEL  
City University of New York  
New York NY, USA  
jterilla@gc.cuny.edu

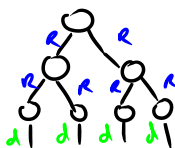
## IV BEYOND TUCKER, TT, ...

TT:

PEPS :  
(2dimTN)



Hierarchical Tucker:  
(Tree TN)



Tensor Ring:

